



Atty. Dkt. No. 016778-0470

Applicant: Masataka ANDOH, et al.  
Title: cDNA MICROARRAY DATA CORRECTION SYSTEM, METHOD,  
PROGRAM, AND MEMORY MEDIUM  
Appl. No.: 10/696,572  
Filing Date: 10/30/2003  
Examiner: Unknown  
Art Unit: 1651

**CLAIM FOR CONVENTION PRIORITY**

Commissioner for Patents  
P.O. Box 1450  
Alexandria, VA 22313-1450

Sir:

The benefit of the filing date of the following prior foreign application filed in the following foreign country is hereby requested, and the right of priority provided in 35 U.S.C. § 119 is hereby claimed.

In support of this claim, filed herewith is a certified copy of said original foreign application:

Japanese Patent Application No. 2003-124585  
filed 04/28/2003.

Respectfully submitted,

Date: July 1, 2004

FOLEY & LARDNER LLP  
Customer Number: 22428  
Telephone: (202) 672-5407  
Facsimile: (202) 672-5399

By

*Philip J. Artisola*

Reg. No.  
38,819

for /

David A. Blumenthal  
Attorney for Applicant  
Registration No. 26,257

US

日 本 国 特 許 庁  
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日                      2 0 0 3 年    4 月 2 8 日  
Date of Application:

出 願 番 号                      特 願 2 0 0 3 - 1 2 4 5 8 5  
Application Number:

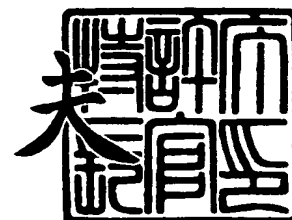
[ST. 10/C] :                      [ J P 2 0 0 3 - 1 2 4 5 8 5 ]

出      願      人  
Applicant(s):                      日本電気株式会社  
   大 瀧   慈  
   社団法人バイオ産業情報化コンソーシアム

2 0 0 3 年    9 月 2 5 日

特許庁長官  
Commissioner,  
Japan Patent Office

今 井 康 夫



出証番号    出証特 2 0 0 3 - 3 0 7 8 9 6 8

【書類名】 特許願

【整理番号】 64002125

【特記事項】 特許法第 3 0 条第 1 項の規定の適用を受けようとする特  
許出願

【提出日】 平成15年 4月28日

【あて先】 特許庁長官殿

【国際特許分類】 G06N 7/00

【発明者】

    【住所又は居所】 東京都港区芝五丁目 7 番 1 号 日本電気株式会社内

    【氏名】 安東 正貴

【発明者】

    【住所又は居所】 東京都港区芝五丁目 7 番 1 号 日本電気株式会社内

    【氏名】 斎藤 彰

【発明者】

    【住所又は居所】 広島県廿日市市宮園 9 丁目 1 の 7

    【氏名】 大瀧 慈

【発明者】

    【住所又は居所】 広島県広島市佐伯区楽々園 5 - 9 五日市住宅 1 6 - 2  
0 2

    【氏名】 佐藤 健一

【発明者】

    【住所又は居所】 広島県広島市南区仁保南 2 丁目 8 - 7

    【氏名】 西山 正彦

【発明者】

    【住所又は居所】 東京都中央区八丁堀二丁目 2 6 番 9 号 グランデビルデ  
ィング 社団法人バイオ産業情報化コンソーシアム内

    【氏名】 大谷 敬子

## 【特許出願人】

【識別番号】 000004237

【氏名又は名称】 日本電気株式会社

## 【特許出願人】

【識別番号】 503077165

【氏名又は名称】 大瀧 慈

## 【特許出願人】

【識別番号】 500535301

【氏名又は名称】 社団法人バイオ産業情報化コンソーシアム

## 【代理人】

【識別番号】 100071272

【弁理士】

【氏名又は名称】 後藤 洋介

## 【選任した代理人】

【識別番号】 100077838

【弁理士】

【氏名又は名称】 池田 憲保

## 【手数料の表示】

【予納台帳番号】 012416

【納付金額】 21,000円

## 【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 0018587

【ブルーフの要否】 要

【書類名】 明細書

【発明の名称】 cDNAマイクロアレイデータの補正システム、方法、プログラム、及び記録媒体

【特許請求の範囲】

【請求項1】 マイクロアレイデータのグローバル及びローカルな歪みに対してより精密な補正をし、さらに蛍光色素の感度の違いによる測定誤差を補正するcDNAマイクロアレイデータの補正システムにおいて、

各スポットのバックグラウンドノイズの除去及び信頼性を示すフラッグ情報を考慮し、あらかじめ調整されている遺伝子発現強度データを入力する入力装置と、前記入力された遺伝子発現強度データに対して、グリッド毎の順序統計量を用いて遺伝子発現強度データを規準化し、規準化された規準化遺伝子発現強度データを送出するデータ規準化手段と、前記規準化遺伝子発現強度データに対して、グリッドの座標におけるスポット位置に依存する歪みをノンパラメトリック平滑化法によって推定し、補正した第1の遺伝子発現強度補正データを送出するスポット位置による第1の補正手段と、前記第1の遺伝子発現強度補正データに対してS-D変換を行い、遺伝子発現強度データに潜在しうる蛍光色素の感度の違いによる歪みをノンパラメトリック平滑化法によって推定し、蛍光色素の感度の違いによる歪みが補正された第2の遺伝子発現強度補正データを送出する第2の補正手段とを具備するデータ解析装置と、

前記第2の遺伝子発現強度補正データを出力する出力装置を

有することを特徴とするcDNAマイクロアレイデータの補正システム。

【請求項2】 さらに、任意の段階で遺伝子発現強度データの歪みを定量化し、S-Dプロット上に視覚化するS-D変換手段を有していることを特徴とする請求項1記載のcDNAマイクロアレイデータの補正システム。

【請求項3】 前記順序統計量は以下の数1（尚、外1は、前記規準化遺伝子発現強度データであり、外2はチャンネルによって得られた全スポットの遺伝子発現強度データであり、 $L_k(c)$ および $M_k(c)$ はそれぞれグリッドkにおいて、チャンネルcによって得られた遺伝子発現強度データの25%点および50%点を示す。）で示されることを特徴とする請求項1又は2記載のcDNAマイ

クロアレイデータの補正システム。

【数 1】

$$w_{ij}^k(c) = \frac{y_{ij}^k(c) - L_k(c)}{M_k(c) - L_k(c)}, \quad c=1, 2, \quad i=1, \dots, I, \quad j=1, \dots, J, \quad k=1, \dots, K.$$

【外 1】

$$w_{ij}^k(c)$$

【外 2】

$$y_{ij}^k(c)$$

【請求項 4】 前記データ規準化手段は、少なくとも 2 つの遺伝子発現強度データチャンネルによって得られた全スポットの遺伝子発現強度データを規準化したかどうかを判定し、全スポットの遺伝子発現強度データを規準化するまで続けることを特徴とする請求項 3 記載の cDNA マイクロアレイデータの補正システム。

【請求項 5】 前記規準化遺伝子発現強度データは、真の発現強度とスポット位置に依存する歪みとの和によって表されることを特徴とする請求項 1 記載の cDNA マイクロアレイデータの補正システム。

【請求項 6】 前記第 1 の補正手段は、スポット位置に依存する歪みを x 軸、y 軸、及び前記 x、y 軸の交互作用による歪み（それぞれ、外 3、外 4、外 5 とする。）の回帰関係で示されるノンパラメトリック回帰モデルにより記述し、以下の数 2 に示されるノンパラメトリック平滑化法を用いて、スポット位置による歪み（外 6）を推定することを特徴とする請求項 1 記載の cDNA マイクロアレイデータの補正システム。

【数 2】

$$\hat{\xi}_{ij}^k(c) = \hat{\alpha}_k^{(c)}(i) + \hat{\beta}_k^{(c)}(j) + \hat{\gamma}_k^{(c)}((i - m_i)(j - m_j)), \quad c=1,2, \quad i=1,\dots,I, \quad j=1,\dots,J.$$

【外 3】

$$\alpha_k^{(c)}(i)$$

【外 4】

$$\beta_k^{(c)}(j)$$

【外 5】

$$\gamma_k^{(c)}((i - m_i)(j - m_j))$$

【外 6】

$$\xi_{ij}^k(c)$$

【請求項 7】 前記スポット位置による歪みの補正は、以下の数 3（尚、外 7 は補正された真の遺伝子発現強度データである。）に従って行なわれることを特徴とする請求項 6 記載の cDNA マイクロアレイデータの補正システム。

【数 3】

$$\hat{z}_{ij}^k(c) = w_{ij}^k(c) - \hat{\xi}_{ij}^k(c)$$

【外 7】

$$\hat{z}_{ij}^k(c)$$

【請求項 8】 前記第 2 の補正手段における前記 S-D 変換は、以下の数 4 に従って行なわれることを特徴とする請求項 7 記載の cDNA マイクロアレイデータの補正システム。

【数 4】

$$u_{ij}^k = \hat{z}_{ij}^k(1) + \hat{z}_{ij}^k(2)$$

$$v_{ij}^k = \hat{z}_{ij}^k(1) - \hat{z}_{ij}^k(2)$$

【請求項 9】 前記第 2 の補正手段は、以下の数 5 で示されるノンパラメトリック回帰モデルにより記述し、以下の数 6 及び数 7 で示されるノンパラメトリック平滑化法を用いて、蛍光色素の感度による測定誤差を推定し、補正を行うことを特徴とする請求項 8 記載の cDNA マイクロアレイデータの補正システム。

【数 5】

$$v_{ij}^k = \phi(u_{ij}^k) + \varepsilon_{ij}^k, \varepsilon_{ij}^k = N(0, v^2)$$

【数 6】

$$\eta_{ij}^k = v_{ij}^k - \hat{\phi}(u_{ij}^k)$$

【数 7】

$$\hat{y}_{ij}^k(1) = \frac{1}{2} (u_{ij}^k + \eta_{ij}^k)$$

$$\hat{y}_{ij}^k(2) = \frac{1}{2} (u_{ij}^k - \eta_{ij}^k)$$

【請求項 10】 前記補正の前提として、遺伝子が発現している確率は 0.5 より小さいと仮定し、各グリッド内の半分以上のスポットで検出される蛍光強度は、バックグラウンドノイズあるいは系統誤差を示しているとする特徴とする請求項 1 記載の cDNA マイクロアレイデータの補正システム。

【請求項 11】 さらに前記補正の前提として、グリッドにおいて、少なくとも 2 つの遺伝子発現強度データチャンネルによって得られた蛍光強度の 25% 点と 50% 点を、それぞれ  $L_k(c)$  および  $M_k(c)$  とするとき、遺伝子の大半は非発現状態にあり全てのグリッドとチャンネルにおいて蛍光強度の 50% 点以下の分布は共通であるという前提に基づき、 $L_k(c)$  と  $M_k(c) - L_k(c)$  は各グリッドおよび各チャンネルで等しいと仮定することを特徴とする請求項 10 記載の cDNA マイクロアレイデータの補正システム。

【請求項 12】 マイクロアレイデータのグローバル及びローカルな歪みに



対してより精密な補正をし、さらに蛍光色素の感度の違いによる測定誤差を補正する cDNA マイクロアレイデータの補正方法において、

各スポットのバックグラウンドノイズの除去及び信頼性を示すフラッグ情報を考慮し、あらかじめ調整されている遺伝子発現強度データを入力するステップと、大半の遺伝子は発現していないことを前提として、前記入力された遺伝子発現強度データに対してグリッド毎の順序統計量を用いて、当該遺伝子発現強度データを規準化するステップと、

前記規準化された規準化遺伝子発現強度データを出力するステップと、

前記規準化遺伝子発現強度データに対して、グリッドの座標におけるスポット位置に依存する歪みをノンパラメトリック平滑化法によって推定し、スポット位置に依存したデータの歪みを補正するステップと、

前記スポット位置に依存したデータ歪みの補正がされた第 1 の遺伝子発現強度補正データを出力するステップと、

前記第 1 の遺伝子発現強度補正データに対して、S-D 変換を行い、遺伝子発現強度データに潜在しうる蛍光色素の感度の違いによる歪みをノンパラメトリック平滑化法によって推定し、蛍光色素の感度の違いによる歪みを補正するステップと、

前記蛍光色素の感度の違いによる歪みの補正がされた第 2 の遺伝子発現強度補正データを出力するステップとを

有することを特徴とする cDNA マイクロアレイデータの補正方法。

【請求項 13】 さらに、任意の段階で遺伝子発現強度データの歪みを定量化し、S-D プロット上に視覚化するステップを有していることを特徴とする請求項 12 記載の cDNA マイクロアレイデータの補正方法。

【請求項 14】 前記順序統計量は以下の数 8（尚、外 8 は、前記規準化遺伝子発現強度データであり、外 9 はチャンネルによって得られた全スポットの遺伝子発現強度データであり、 $L_k(c)$  および  $M_k(c)$  はそれぞれグリッド  $k$  において、チャンネル  $c$  によって得られた遺伝子発現強度データの 25% 点および 50% 点を示す。）で示されることを特徴とする請求項 12 又は 13 記載の cDNA マイクロアレイデータの補正方法。

【数 8】

$$w_{ij}^k(c) = \frac{y_{ij}^k(c) - L_k(c)}{M_k(c) - L_k(c)}, \quad c=1, 2, \quad i=1, \dots, I, \quad j=1, \dots, J, \quad k=1, \dots, K.$$

【外 8】

$$w_{ij}^k(c)$$

【外 9】

$$y_{ij}^k(c)$$

【請求項 15】 前記データを規準化するステップにおいて、少なくとも 2 つの遺伝子発現強度データチャンネルによって得られた全スポットの遺伝子発現強度データを規準化したかどうかを判定し、全スポットの遺伝子発現強度データを規準化するまで続けることを特徴とする請求項 14 記載の cDNA マイクロアレイデータの補正方法。

【請求項 16】 前記規準化遺伝子発現強度データは、真の発現強度とスポット位置に依存する歪みとの和によって表されることを特徴とする請求項 15 記載の cDNA マイクロアレイデータの補正方法。

【請求項 17】 前記スポット位置に依存したデータの歪みを補正するステップにおいて、スポット位置に依存する歪みを x 軸、y 軸、及び前記 x, y 軸の交互作用による歪み（それぞれ、外 10, 外 11, 外 12 とする。）の回帰関係で示されるノンパラメトリック回帰モデルにより記述し、以下の数 9 で示されるノンパラメトリック平滑化法を用いて、スポット位置による歪み（外 13）を推定することを特徴とする請求項 12 記載の cDNA マイクロアレイデータの補正方法。

【数 9】

$$\hat{\xi}_{ij}^k(c) = \hat{\alpha}_k^{(c)}(i) + \hat{\beta}_k^{(c)}(j) + \hat{\gamma}_k^{(c)}((i - m_i)(j - m_j)), \quad c=1,2, i=1,\dots,I, j=1,\dots,J.$$

【外 1 0】

$$\alpha_k^{(c)}(i)$$

【外 1 1】

$$\beta_k^{(c)}(j)$$

【外 1 2】

$$\gamma_k^{(c)}((i - m_i)(j - m_j))$$

【外 1 3】

$$\xi_{ij}^k(c)$$

【請求項 1 8】 前記スポット位置による歪みの補正は、以下の数 1 0（尚、外 1 4 は補正された真の遺伝子発現強度データである。）に従って行なわれることを特徴とする請求項 1 7 記載の c DNA マイクロアレイデータの補正方法。

【数 1 0】

$$\hat{z}_{ij}^k(c) = w_{ij}^k(c) - \hat{\xi}_{ij}^k(c)$$

【外 1 4】

$$\hat{z}_{ij}^k(c)$$

【請求項 1 9】 前記蛍光色素の感度の違いによる歪みを補正するステップにおける前記 S-D 変換は、以下の数 1 1 に従って行なわれることを特徴とする請求項 1 8 記載の c DNA マイクロアレイデータの補正方法。

【数 1 1】

$$u_{ij}^k = \hat{z}_{ij}^k(1) + \hat{z}_{ij}^k(2)$$

$$v_{ij}^k = \hat{z}_{ij}^k(1) - \hat{z}_{ij}^k(2)$$

【請求項 2 0】 前記蛍光色素の感度の違いによる歪みを補正するステップにおいて、以下の数 1 2 で示されるノンパラメトリック回帰モデルにより記述し、以下の数 1 3 及び数 1 4 で示されるノンパラメトリック平滑化法を用いて、蛍光色素の感度による測定誤差を推定し、補正を行うことを特徴とする請求項 1 9 記載の c DNA マイクロアレイデータの補正方法。

【数 1 2】

$$v_{ij}^k = \phi(u_{ij}^k) + \varepsilon_{ij}^k, \varepsilon_{ij}^k = N(0, v^2)$$

【数 1 3】

$$\eta_{ij}^k = v_{ij}^k - \hat{\phi}(u_{ij}^k)$$

【数 1 4】

$$\hat{y}_{ij}^k(1) = \frac{1}{2} (u_{ij}^k + \eta_{ij}^k)$$

$$\hat{y}_{ij}^k(2) = \frac{1}{2} (u_{ij}^k - \eta_{ij}^k)$$

【請求項 2 1】 前記補正の前提として、遺伝子が発現している確率は 0.5 より小さいと仮定し、各グリッド内の半分以上のスポットで検出される蛍光強度は、バックグラウンドノイズあるいは系統誤差を示しているとする特徴とする請求項 1 2 記載の c DNA マイクロアレイデータの補正方法。

【請求項 2 2】 さらに前記補正の前提として、グリッドにおいて、少なくとも 2 つの遺伝子発現強度データチャンネルによって得られた蛍光強度の 25 % 点と 50 % 点を、それぞれ  $L_k(c)$  および  $M_k(c)$  とするとき、遺伝子の大半は非発現状態にあり、全てのグリッドとチャンネルにおいて蛍光強度の 50 % 点以下の分布は共通であるという前提に基づき、 $L_k(c)$  と  $M_k(c) - L_k(c)$  は各グリッドおよび各チャンネルで等しいと仮定することを特徴とする請求項 2 1 記載の c DNA マイクロアレイデータの補正方法。

【請求項 2 3】 マイクロアレイデータのグローバル及びローカルな歪みに対してより精密な補正をし、さらに蛍光色素の感度の違いによる測定誤差を補正するためコンピュータに、

各スポットのバックグラウンドノイズの除去及び信頼性を示すフラッグ情報を考慮し、あらかじめ調整されている遺伝子発現強度データを入力するステップと、  
大半の遺伝子は発現していないことを前提として、前記入力された遺伝子発現強度データに対してグリッド毎の順序統計量を用いて、当該遺伝子発現強度データを規準化するステップと、

前記規準化された規準化遺伝子発現強度データを出力するステップと、

前記規準化遺伝子発現強度データに対して、グリッドの座標におけるスポット位置に依存する歪みをノンパラメトリック平滑化法によって推定し、スポット位置に依存したデータの歪みを補正するステップと、

前記スポット位置に依存したデータ歪みの補正がされた第 1 の遺伝子発現強度補正データを出力するステップと、

前記第 1 の遺伝子発現強度補正データに対して、S-D 変換を行い、遺伝子発現強度データに潜在しうる蛍光色素の感度の違いによる歪みをノンパラメトリック平滑化法によって推定し、蛍光色素の感度の違いによる歪みを補正するステップと、

前記蛍光色素の感度の違いによる歪みの補正がされた第 2 の遺伝子発現強度補正データを出力するステップとを

を実行させるための cDNA マイクロアレイデータ補正プログラム。

【請求項 2 4】 マイクロアレイデータのグローバル及びローカルな歪みに対してより精密な補正をし、さらに蛍光色素の感度の違いによる測定誤差を補正するためコンピュータに、

各スポットのバックグラウンドノイズの除去及び信頼性を示すフラッグ情報を考慮し、あらかじめ調整されている遺伝子発現強度データを入力するステップと、  
大半の遺伝子は発現していないことを前提として、前記入力された遺伝子発現強度データに対してグリッド毎の順序統計量を用いて、当該遺伝子発現強度データを規準化するステップと、

前記規準化された規準化遺伝子発現強度データを出力するステップと、

前記規準化遺伝子発現強度データに対して、グリッドの座標におけるスポット位置に依存する歪みをノンパラメトリック平滑化法によって推定し、スポット位置に依存したデータの歪みを補正するステップと、

前記スポット位置に依存したデータ歪みの補正がされた第 1 の遺伝子発現強度補正データを出力するステップと、

前記第 1 の遺伝子発現強度補正データに対して、S-D変換を行い、遺伝子発現強度データに潜在しうる蛍光色素の感度の違いによる歪みをノンパラメトリック平滑化法によって推定し、蛍光色素の感度の違いによる歪みを補正するステップと、

前記蛍光色素の感度の違いによる歪みの補正がされた第 2 の遺伝子発現強度補正データを出力するステップとを

を実行させるための cDNA マイクロアレイデータ補正プログラムを記録したコンピュータ読み取り可能な記録媒体。

#### 【発明の詳細な説明】

##### 【0001】

#### 【発明の属する技術分野】

本発明は、数理モデルに基づいた cDNA マイクロアレイデータのデータ補正システム、方法、プログラム及び記録媒体に関し、特にグローバルノーマライゼーションとローカルノーマライゼーション、さらに蛍光色素の感度の違いによる測定の変形の補正をすることができる cDNA マイクロアレイデータの補正システム、方法、プログラム及び記録媒体に関するものである。

##### 【0002】

#### 【従来の技術】

現在、ゲノム研究は個々の遺伝子についての構造解析から体系的な遺伝子の機能解析へと展開しつつある。機能未知の遺伝子や総体としての遺伝子の機能解析のために、多数の遺伝子の発現強度を同時に定量化することのできる cDNA (相補的な DNA) マイクロアレイを用いた実験はその有効性が大いに期待されている。

**【0003】**

二色蛍光法による cDNA マイクロアレイを用いた実験の目的は二種類の細胞の遺伝子発現の違いを検出することにある。ここで、二色蛍光法による cDNA マイクロアレイの概要について述べる。まず、多数の遺伝子セットの cDNA を参照用のプローブとして、スライドガラス上にアレイ状に高密度に固定化する（マイクロアレイ）。

**【0004】**

次に、条件の異なる 2 種類のサンプル、細胞 1 と細胞 2（例えば正常細胞と癌細胞）から抽出した mRNA をそれぞれ波長の異なる蛍光色素でラベルし、ターゲット cDNA を合成する。そして、それらを等量混合したものをマイクロアレイに固定化された参照用のプローブ cDNA に競合的にハイブリダイズさせる。ハイブリダイゼーション後、スキャナーでそれぞれの蛍光色素強度を測定する。細胞 1 にラベルされた蛍光色素をチャンネル 1 により、細胞 2 にラベルされた蛍光色素をチャンネル 2 により読み取り、それぞれを各細胞の遺伝子発現強度データ（マイクロアレイデータ）とする。

**【0005】**

このように、マイクロアレイデータが得られるまでの過程は複雑であり、高度な実験技術が必要とされることから、実験の各段階において様々な実験誤差が生じると考えられる。このため、マイクロアレイデータから真に生物学的意味のあるデータを取り出すためには遺伝子発現強度の分布と実験誤差の解析は解決すべき重要な課題である。

**【0006】**

遺伝子発現強度の分布に関しては、例えば、以下の非特許文献 1 を参照すると、Newton 等は遺伝子発現強度にガンマ分布関数を仮定し、遺伝子発現強度比（チャンネル 1 とチャンネル 2 の遺伝子発現強度データの比）についての統計学的性質を考察している。

**【0007】**

また、観測された遺伝子発現強度データに対しては、例えば、以下の非特許文献 2 を参照すると、Lee 等は真の遺伝子発現強度を 2 個の水準値に分離できる

ことおよび偶然誤差の存在を前提として、以下の数 1 5 に示されるような混合正規分布を適用し、遺伝子発現強度データについての統計学的考察を行った。

【0 0 0 8】

【数 1 5】

$$f(x) = p\phi(x - \mu_1 | \sigma_1^2) + (1-p)\phi(x - \mu_2 | \sigma_2^2)$$

ここで、 $x$  はスキャナーなどによって得られる蛍光強度などの遺伝子発現強度データを表し、右辺第 1 項の外 1 5 は

【外 1 5】

$$\phi(x - \mu_1 | \sigma_1^2)$$

遺伝子が発現しているときの平均  $\mu_1$ 、分散外 1 6 の正規分布、

【外 1 6】

$$\sigma_1^2$$

また、同第 2 項の外 1 7 は遺伝子が発現していないときの平均  $\mu_2$ 、分散外 1 8 の正規分布の密度関数を表し、

【外 1 7】

$$\phi(x - \mu_2 | \sigma_2^2)$$

【外 1 8】

$$\sigma_2^2$$

$p$  はその混合率を表す母数である。

【0 0 0 9】

実験誤差の解析については、系統誤差の除去、いわゆるノーマライゼーション



の方法がいくつか提案されている。ノーマライゼーションの方法は、大きく分けてアレイ上のすべてのスポットを対象にしたグローバルノーマライゼーションと、あるサブセットに分けた（例えばグリッド単位の）スポットを対象にしたローカルノーマライゼーションの二つが提案されている。グローバルノーマライゼーションについては、例えば、以下の非特許文献 3 を参照すると、C h e n 等は二つの細胞の遺伝子発現強度の中央値は等しいとしてチャンネル 1 とチャンネル 2 で得られた測定値の補正を行った。ローカルノーマライゼーションについては、例えば、以下の非特許文献 4、5、6 を参照すると、D u d o i t や S c h u c h h a r d t や Y a n g は、系統誤差が、スポットのスライドガラス上の位置や、二種類の蛍光色素の感度の違いによって生じたものと考え、それらを除去する方法を提案した。

**【 0 0 1 0 】****【非特許文献 1】**

Newton et. al、2001年、ジャーナル・オブ・コンピュテーショナル・バイオロジー、第8巻、37～52頁（Journal of Computational Biology Vol. 8, pp. 37-52）

**【 0 0 1 1 】****【非特許文献 2】**

Lee et. al、2000年、プロシーディング・オブ・ザ・ナショナル・アカデミー・オブ・サイエンス、第97巻、第18号、9834～9839頁（Proceeding of the National Academy of Sciences Vol. 97, No 18, pp. 9834-9839）

**【 0 0 1 2 】****【非特許文献 3】**

Chen et. al、1997年、ジャーナル・オブ・バイオメディカル・オプティクス、第 2 号、364～374頁（Journal of Biomedical Optics Vol. 2, pp. 364-374）

**【 0 0 1 3 】****【非特許文献 4】**

Dudoit et. al、2000. Statistical methods for identifying differentia

lly expressed genes in replicated cDNA microarray experiments. Technical  
~Report #578 2.

【 0 0 1 4 】

【非特許文献 5】

Schuchhardt et. al、2000年、ヌクレ・アシッド・リサーチ、第28巻、第10  
号 (Nucleic Acids Research, 2000, Vol.28, No. 10)

【 0 0 1 5 】

【非特許文献 6】

Yang et. al、2002年、ヌクレ・アシッド・リサーチ、第30巻、第4号 (Nuc  
leic Acids Research, 2002, Vol.30, No. 4)

【 0 0 1 6 】

【発明が解決しようとする課題】

上記した従来技術における問題点は、マイクロアレイデータの解析結果は再現  
性に乏しく不安定なものになりがちで、精度や効率は低いものとみなされている  
ことである。その理由は、遺伝子の発現に関する真の信号と実験誤差の分離が十  
分に行われていないからである。その背景要因として、遺伝子の発現強度はそれ  
ぞれの遺伝子によってレベルが異なっていることが考えられ、その場合、上記数  
1 5 によるモデルは明らかに過大に単純化されすぎている。

【 0 0 1 7 】

本発明の目的は、マイクロアレイ上の遺伝子発現強度データに関してよりもっ  
ともな数理モデルを想定して、グローバルおよびローカルな歪みに対して高い精  
度の補正を行い、さらに蛍光色素の感度の違いによる測定誤差を補正するための  
包括的なノーマライゼーションの方法およびシステムを提供することである。

【 0 0 1 8 】

【課題を解決するための手段】

本発明の c DNA マイクロアレイデータの補正システムは、蛍光強度などの遺  
伝子発現強度データを入力する入力装置と、プログラムの制御により動作するデ  
ータ解析装置と、出力装置とを含む。なお、入力される遺伝子発現強度データは  
、各スポットのバックグラウンドノイズの除去や各スポットの信頼性を示すフラ

ッグ情報を考慮し、あらかじめ調整されているものとする。

#### 【0019】

前記データ解析装置は、下記の三個の連続した処理過程で構成される。第一処理過程であるデータ規準化手段では、前記入力装置から遺伝子発現強度データを入力し、大半の遺伝子は発現していないことを前提としてグリッド毎の順序統計量を用いて遺伝子発現強度データを規準化し、規準化した遺伝子発現強度データを出力する。

#### 【0020】

第二処理過程であるスポット位置による補正手段では、前記規準化された遺伝子発現強度データを入力し、グリッド毎にスポット位置に依存する歪みをノンパラメトリック平滑化法によって推定し、スポット位置に依存したデータの歪みを補正した遺伝子発現強度データを出力する。

#### 【0021】

第三処理過程であるS-Dプロットによる補正手段では、第二処理過程の段階まで補正された遺伝子発現強度データに対してMA変換の変形であるS-D変換(MA変換およびMAプロットについては、上記非特許文献6を参照)を行い、遺伝子発現強度データに潜在しうる蛍光色素の感度の違いによる歪みをノンパラメトリック平滑化法によって推定し、蛍光色素の感度の違いによる歪みを補正した遺伝子発現強度データを前記出力装置に出力する。

#### 【0022】

なお、本システムは、任意の段階で遺伝子発現強度データの歪みを定量化し、S-Dプロット上に視覚化するS-D変換手段を有していることを特徴とする。

#### 【0023】

このような構成を採用し、遺伝子発現強度データを補正することにより、本発明の目的を達成することができる。

#### 【0024】

##### 【発明の実施の形態】

はじめに、本発明におけるマイクロアレイの構造を説明する。図1を参照すると、K個の各グリッドに $I \times J$ 個ずつ、合計 $K \times I \times J$ 個のcDNAがスライド

ガラス上にスポットされている。いま、グリッド  $k$  における座標  $(i, j)$  にスポットされた  $c$  DNA に対して、チャンネル  $c = 1, 2$  によって得られた蛍光強度を外 19 とする。

【0025】

【外 19】

$$y_{ij}^k(c), c = 1, 2$$

次に、以下の 2 つの仮定をする。

【0026】

(仮定 1)

遺伝子が発現している確率は 0.5 より小さいと仮定し、各グリッド内の半分以上のスポットで検出される蛍光強度外 20 は、バックグラウンドノイズあるいは系統誤差を示しているとする。

【0027】

【外 20】

$$y_{ij}^k(c)$$

(仮定 2)

グリッド  $k$  において、チャンネル  $c$  によって得られた蛍光強度外 21 の 25 % 点と 50 % 点を、

【外 21】

$$y_{ij}^k(c)$$

それぞれ  $L_k(c)$  および  $M_k(c)$  とするとき、遺伝子の大半は非発現状態にあり全てのグリッドとチャンネルにおいて蛍光強度の 50 % 点以下の分布は共通であるという前提に基づき、 $L_k(c)$  と  $M_k(c) - L_k(c)$  は各グリッドおよび各チャンネルで等しいと仮定する。

**【0 0 2 8】**

次に、以上の仮定をもとに、本発明の第 1 の実施の形態について図面を参照して詳細に説明する。図 2 を参照すると、本発明の第 1 の実施の形態は、蛍光強度などの遺伝子発現強度データを入力する入力装置 1 と、プログラム制御により動作するデータ解析装置 2 と、ディスプレイ装置や印刷装置等の出力装置 3 とを含む。データ解析装置は、データ規準化手段 2 1 と、スポット位置による補正手段 2 2 と、S-D プロットによる補正手段 2 3 とを備えている。

**【0 0 2 9】**

データ規準化手段 2 1 は、与えられた遺伝子発現強度データに対して、グリッド毎の順序統計量を用いて遺伝子発現強度データを規準化し、スポット位置による補正手段 2 2 及び S-D 変換手段 2 4 に送る。

**【0 0 3 0】**

スポット位置による補正手段 2 2 は、データ規準化手段 2 1 から送られてきた規準化された遺伝子発現強度データに対して、グリッド毎にスポット位置に依存する歪みをノンパラメトリック平滑化法によって推定し、補正した遺伝子発現強度データを S-D プロットによる補正手段 2 3 及び S-D 変換手段 2 4 に送る。

**【0 0 3 1】**

S-D プロットによる補正手段 2 3 は、スポット位置による補正手段 2 2 から送られてきた補正された遺伝子発現強度データに S-D 変換を行い、蛍光色素の感度の違いに起因する歪みをノンパラメトリック平滑化法により補正した後、遺伝子発現強度データを出力装置 3 へ送る。

**【0 0 3 2】**

S-D 変換手段 2 4 は送られてきた遺伝子発現強度データに S-D 変換を行い、出力装置 3 へ送る。

**【0 0 3 3】**

次に、図 2、図 3 を参照して本実施の形態について詳細に説明する。入力装置 1 より入力された蛍光強度などの遺伝子発現強度データはデータ規準化手段 2 1 へ送られる。データ規準化手段 2 1 は、送られてきた発現強度データに対して、以下の数 1 6 で示されるように、グリッド毎の順序統計量を用いて発現強度デー

タを規準化する（図3のステップA1）。

【0034】

【数16】

$$w_{ij}^k(c) = \frac{y_{ij}^k(c) - L_k(c)}{M_k(c) - L_k(c)}, \quad c=1,2, i=1,\dots,I, j=1,\dots,J, k=1,\dots,K.$$

2つのチャンネルによって得られた全スポットの遺伝子発現強度データ外22を規準化したかどうかを判定し、

【外22】

$$y_{ij}^k(c)$$

全スポットの遺伝子発現強度データ（ $2 \times I \times J \times K$ 個）を規準化するまで続ける（ステップA2）。

【0035】

データ規準化手段21において規準化された遺伝子発現強度データ外23に対して、

【外23】

$$w_{ij}^k(c)$$

外24を真の発現強度を反映した蛍光強度（以下、真の発現蛍光強度）とし、

【外24】

$$z_{ij}^k(c)$$

外25をグリッドkの座標（i, j）におけるスポット位置に依存する歪みとする。

【0036】

【外 2 5】

$$\xi_{ij}^k(c)$$

このとき、以下の数 1 7 に示すように、遺伝子発現強度データ外 2 6 は、真の発現強度外 2 7 とスポット位置に依存する歪み外 2 8 との和によって表されるとする。

【0 0 3 7】

【外 2 6】

$$w_{ij}^k(c)$$

【外 2 7】

$$z_{ij}^k(c)$$

【外 2 8】

$$\xi_{ij}^k(c)$$

【数 1 7】

$$w_{ij}^k(c) = z_{ij}^k(c) + \xi_{ij}^k(c) + \varepsilon_{ij}^k(c), \varepsilon_{ij}^k(c) \sim N(0, \sigma_k(c)^2), c = 1, 2.$$

ただし、外 2 9 はランダムなノイズであるとする。

【0 0 3 8】

【外 2 9】

$$\varepsilon_{ij}^k(c)$$

スポット位置による補正手段 2 2 は、以下の数 1 8 に示すようにスポット位置

に依存する歪み外 3 0 を「x 軸」、「y 軸」および「2 つの軸の交互作用」による歪みの回帰関係で示されるノンパラメトリック回帰モデルにより記述し、

【外 3 0】

$$\xi_{ij}^k(c)$$

以下の数 1 9 に示すようにノンパラメトリック平滑化法を用いて、スポット位置による歪み外 3 1 を推定する。

【0 0 3 9】

【外 3 1】

$$\xi_{ij}^k(c)$$

【数 1 8】

$$\xi_{ij}^k(c) = \alpha_k^{(c)}(i) + \beta_k^{(c)}(j) + \gamma_k^{(c)}((i - m_i)(j - m_j)), \quad c=1,2, \quad i=1, \dots, I, \quad j=1, \dots, J,$$

$$\sum_i \alpha_k^{(c)}(i) = 0, \quad \sum_j \beta_k^{(c)}(j) = 0, \quad \sum_i \sum_j \gamma_k^{(c)}((i - m_i)(j - m_j)) = 0.$$

【数 1 9】

$$\hat{\xi}_{ij}^k(c) = \hat{\alpha}_k^{(c)}(i) + \hat{\beta}_k^{(c)}(j) + \hat{\gamma}_k^{(c)}((i - m_i)(j - m_j)), \quad c=1,2, \quad i=1, \dots, I, \quad j=1, \dots, J.$$

ここで、外 3 2 とする。

【0 0 4 0】

【外 3 2】

$$m_i = \lfloor I / 2 \rfloor, \quad m_j = \lfloor J / 2 \rfloor$$

外 3 3 は  $\alpha$  以上の最小の整数とする。



【0041】

【外33】

[ $\alpha$ ]

スポット位置による補正手段22は、以下の数20に示すように、データ規準化手段21において規準化された遺伝子発現強度データ外34に対して、推定されたスポット位置による歪み外35を補正する（ステップA3）。

【0042】

【外34】

 $w_{ij}^k(c)$ 

【外35】

 $\hat{\xi}_{ij}^k(c)$ 

【数20】

$$\hat{z}_{ij}^k(c) = w_{ij}^k(c) - \hat{\xi}_{ij}^k(c)$$

データ規準化手段21において規準化された全スポットの遺伝子発現強度データ外36に対して、

【外36】

 $w_{ij}^k(c)$ 

スポット位置による歪み外37の補正をしたかどうかを判定し、

【外37】

 $\hat{\xi}_{ij}^k(c)$ 

全スポットの遺伝子発現強度データ（ $2 \times I \times J \times K$ 個）を補正するまで続ける（ステップA4）。

【0043】

S-Dプロットによる補正手段 2 3 は、スポット位置による補正手段 2 2 において補正された真の遺伝子発現強度データ外 3 8 に対して、

【外 3 8】

$$\hat{z}_{ij}^k(c)$$

以下の数 2 1 に示すように、S-D変換を行う。

【0 0 4 4】

【数 2 1】

$$u_{ij}^k = \hat{z}_{ij}^k(1) + \hat{z}_{ij}^k(2)$$

$$v_{ij}^k = \hat{z}_{ij}^k(1) - \hat{z}_{ij}^k(2)$$

さらに、以下の数 2 2 で示されるようなノンパラメトリック回帰モデルを記述し、以下の数 2 3 及び数 2 4 に示すようにノンパラメトリック平滑化法を用いて蛍光色素の感度による測定誤差を推定し、補正を行う（ステップ A 5）。

【0 0 4 5】

【数 2 2】

$$v_{ij}^k = \phi(u_{ij}^k) + \varepsilon_{ij}^k, \varepsilon_{ij}^k = N(0, v^2)$$

【数 2 3】

$$\eta_{ij}^k = v_{ij}^k - \hat{\phi}(u_{ij}^k)$$

【数 2 4】

$$\hat{y}_{ij}^k(1) = \frac{1}{2} (u_{ij}^k + \eta_{ij}^k)$$

$$\hat{y}_{ij}^k(2) = \frac{1}{2} (u_{ij}^k - \eta_{ij}^k)$$

スポット位置による補正手段 2 2 において補正された真の遺伝子発現強度データ外 3 9 に対して、

【外 3 9】

$$\hat{z}_{ij}^k(c)$$

S-Dプロットによる補正をしたかどうかを判定し、全スポットの真の遺伝子発現強度データ ( $2 \times I \times J \times K$  個) を補正するまで続ける (ステップ A 6)。

#### 【0046】

なお、図3のA2、A4の各ステップ終了後、遺伝子発現強度データはS-D変換手段24を介して出力装置3に送られ、S-Dプロットによって遺伝子発現強度データの歪みを視覚化することができる。

#### 【0047】

次に、本実施の形態の効果について説明する。本実施の形態では、グリッド間での順序統計量を用いた規準化 (グローバルノーマライゼーション) とグリッド内でのスポット位置に依存する歪みの補正 (ローカルノーマライゼーション) を組み合わせたノーマライゼーションを行った。これにより、グリッド間での遺伝子発現強度の偏りによる系統誤差と、グリッド内でのスポット位置に依存する歪みを同時に補正することができる。さらに、S-Dプロットによる補正においては、発現強度データの和と差を用いることにより、蛍光色素の感度の違いによる測定誤差を補正することができる。

#### 【0048】

次に、本発明の第2の実施の形態について図面を参照して詳細に説明する。図4を参照すると、本発明の第2の実施の形態は、本発明の第1の実施の形態と同様に、入力装置、データ解析装置、出力装置を備え、更に、データ解析プログラムを記録した記録媒体4を備える。この記録媒体4は可搬形あるいは固定型のいずれであってもよく、磁気ディスク、半導体メモリ、CD-ROMその他の記録媒体であってもよい。

#### 【0049】

また、本手法を実行できるコンピュータプログラムを、ネットワークに接続されたコンピュータの記録装置に格納しておき、ネットワークを介して他のコンピュータに転送することもできる。本アルゴリズムを実行するコンピュータプログラムを提供する提供媒体としては、様々な形式のコンピュータに読み出し可能な媒体として頒布可能であって、特定のタイプの媒体に限定されるものではない。データ解析プログラムは記録媒体4からデータ解析装置5に読み込まれ、データ

解析装置 2 の動作を制御し、入力装置 1 から入力されたデータファイルに対して第 1 の実施の形態におけるデータ処理装置 2 による処理と同一の処理を実行する。

#### 【0050】

##### 【実施例】

以下、本発明の実施例について説明する。例として用いたデータは、異なる 2 種類の癌細胞（A 細胞、B 細胞）の遺伝子発現状況の比較のために行われた実験から得られたものである。

#### 【0051】

一枚のチップ上に 48 グリッド、1 グリッドあたり 441（ $21 \times 21$ ）スポット、計 21168 の遺伝子の発現パターンについて調べたものである。

#### 【0052】

図 5、図 7 はチャンネル 1 により得られたオリジナルデータの A 細胞遺伝子発現強度を示し、図 6、図 8 はチャンネル 2 によって得られたオリジナルデータの B 細胞遺伝子発現強度を示す。それぞれの図は、マイクロアレイ上のスポット位置に対する遺伝子発現強度の対数値をプロットしたものである。また、図 7、図 8 は第 1 グリッドから第 4 グリッドまでを拡大したものである。図 5～図 8 を見ると、遺伝子発現強度がグリッドごとに周期的に繰り返される系統的な歪みが観察される。マイクロアレイ上の遺伝子は無作為にスポットされているので、このような歪みは実験誤差と考えられる。

#### 【0053】

図 9 は、それらの S-D プロットである。横軸は、各チャンネルの遺伝子発現強度の和、縦軸はそれらの差をとったものを示している。各チャンネルの遺伝子発現強度の和が小さい領域と大きい領域においては、各チャンネルの遺伝子発現強度の差は真の遺伝子発現の違いによる影響は小さく、各チャンネルの蛍光色素の感度の違いによるものと考えられる。これにより、図 9 において蛍光色素の感度の違いによって生じたと考えられる歪みが観察される。

#### 【0054】

図 10 に、チャンネル 1 におけるオリジナルデータのスポット位置に対する遺

伝子発現強度の図を示す。図 1 1 に、チャンネル 1 における第一処理過程後のスポット位置に対する遺伝子発現強度の図を示す。図 1 2 に、チャンネル 1 における第二処理過程後のスポット位置に対する遺伝子発現強度の図を示す。スポット位置に依存していたグリッドごとに周期的に繰り返される系統的な歪みが補正されて取り除かれていることがわかる。

#### 【0 0 5 5】

図 1 3 にチャンネル 1 における第三処理過程後のスポット位置に対する遺伝子発現強度の図を示す。図 1 4 ～図 1 7 にチャンネル 2 におけるオリジナルデータ、第一処理過程後、第二処理過程後、第三処理過程後のスポット位置に対する遺伝子発現強度の図を示す。チャンネル 1 と同様にスポット位置に依存していたグリッドごとに周期的に繰り返される系統的な歪みが補正されて取り除かれていることがわかる。

#### 【0 0 5 6】

図 1 8 ～図 2 1 にオリジナルデータ、第一処理過程後、第二処理過程後、第三処理過程後の S - D プロットを示す。図 2 1 を見ると、蛍光色素の感度の違いによる歪みが補正されて取り除かれていることがわかる。

#### 【0 0 5 7】

##### 【発明の効果】

本発明によれば、グリッド間での位置および尺度の揺らぎに対する頑健な順序統計量の 2 5 % 点と 5 0 % 点による規準化（グローバルノーマライゼーション）と、グリッド内でのスポット位置に依存する歪みの補正（ローカルノーマライゼーション）を組み合わせることでノーマライゼーションを行っているため、グリッド間での遺伝子発現強度の偏りや感度の揺らぎによる系統誤差と、グリッド内でのスポット位置に依存する歪みを、発現している遺伝子の頻度や外れ値の影響をほとんど受けることなく同時に補正することができる。

#### 【0 0 5 8】

又、本発明によれば、S - D プロットにおいて遺伝子発現強度データの和と差を用いることによって、それぞれの蛍光色素の感度の違いが得られ易く、それによる測定誤差を的確に抽出することができるため、蛍光色素の感度の違いによる

測定の変みを効率良く補正することができる。

【図面の簡単な説明】

【図 1】

本発明におけるマイクロアレイの構造を示す図である。

【図 2】

本発明の第 1 の実施の形態の構成を示すブロック図である。

【図 3】

本発明の第 1 の実施の形態の動作を示す流れ図である。

【図 4】

本発明の第 2 の実施の形態の構成を示すブロック図である。

【図 5】

チャンネル 1 で得られたオリジナルデータの遺伝子発現強度の図である。

【図 6】

チャンネル 2 で得られたオリジナルデータの遺伝子発現強度の図である。

【図 7】

チャンネル 1 で得られたオリジナルデータ（第 1 グリッドから第 4 グリッド）の遺伝子発現強度の図である。

【図 8】

チャンネル 2 で得られたオリジナルデータ（第 1 グリッドから第 4 グリッド）の遺伝子発現強度の図である。

【図 9】

オリジナルデータに対する S-D プロットである。

【図 10】

チャンネル 1 のオリジナルデータの遺伝子発現強度の図である。

【図 11】

チャンネル 1 の第一処理過程後の遺伝子発現強度の図である。

【図 12】

チャンネル 1 の第二処理過程後の遺伝子発現強度の図である。

【図 13】

チャンネル 1 の第三処理過程後の遺伝子発現強度の図である。

【図 1 4】

チャンネル 2 のオリジナルデータの遺伝子発現強度の図である。

【図 1 5】

チャンネル 2 の第一処理過程後の遺伝子発現強度の図である。

【図 1 6】

チャンネル 2 の第二処理過程後の遺伝子発現強度の図である。

【図 1 7】

チャンネル 2 の第三処理過程後の遺伝子発現強度の図である。

【図 1 8】

オリジナルデータに対する S - D プロットである。

【図 1 9】

第一処理過程後の S - D プロットである。

【図 2 0】

第二処理過程後の S - D プロットである。

【図 2 1】

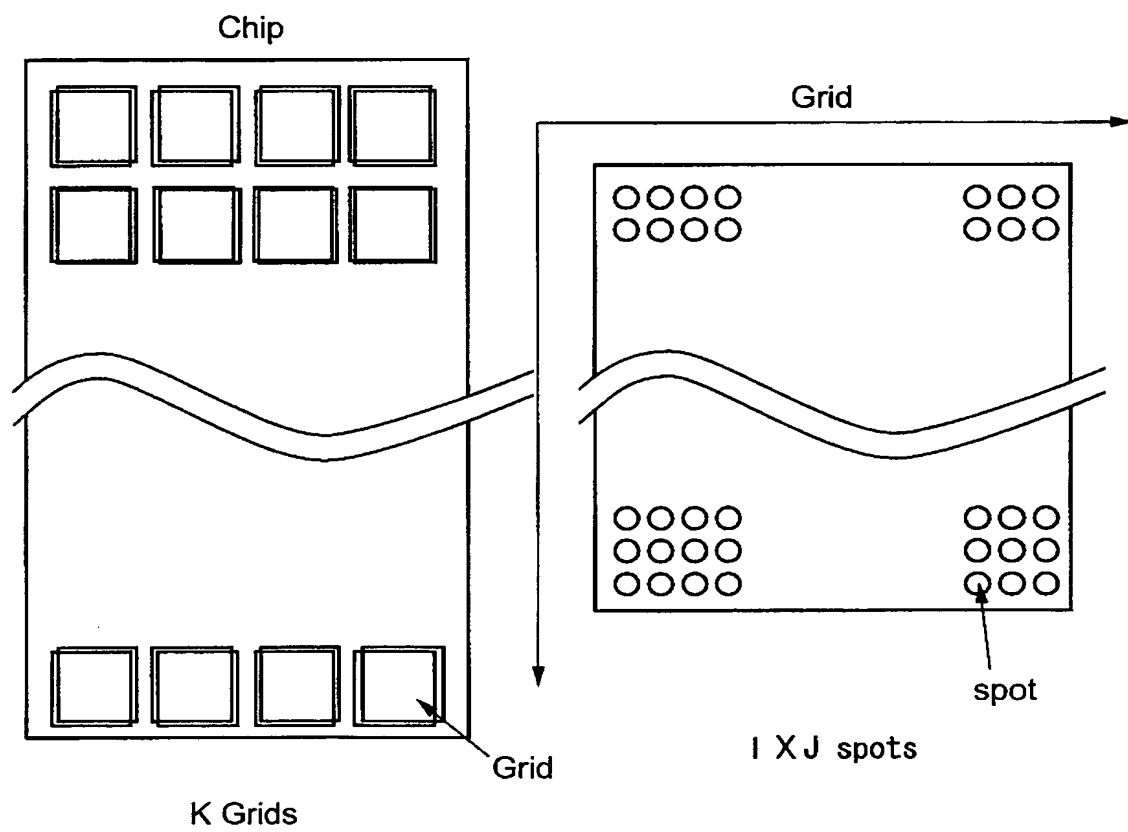
第三処理過程後の S - D プロットである。

【符号の説明】

- 1 入力装置
- 2 データ解析装置
- 3 出力装置
- 2 1 データ規準化手段
- 2 2 スポット位置による補正手段
- 2 3 S - D プロットによる補正手段
- 2 4 S - D 変換手段

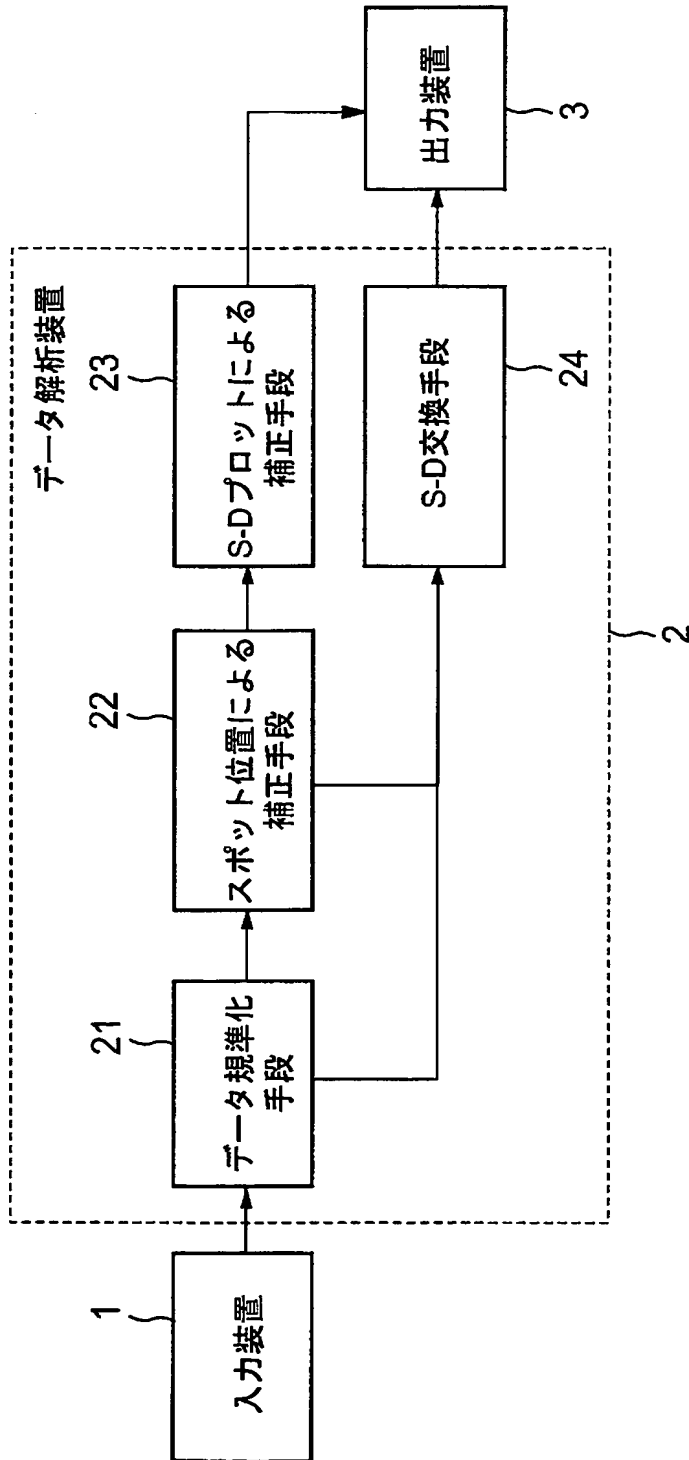
【書類名】 図面

【図 1】

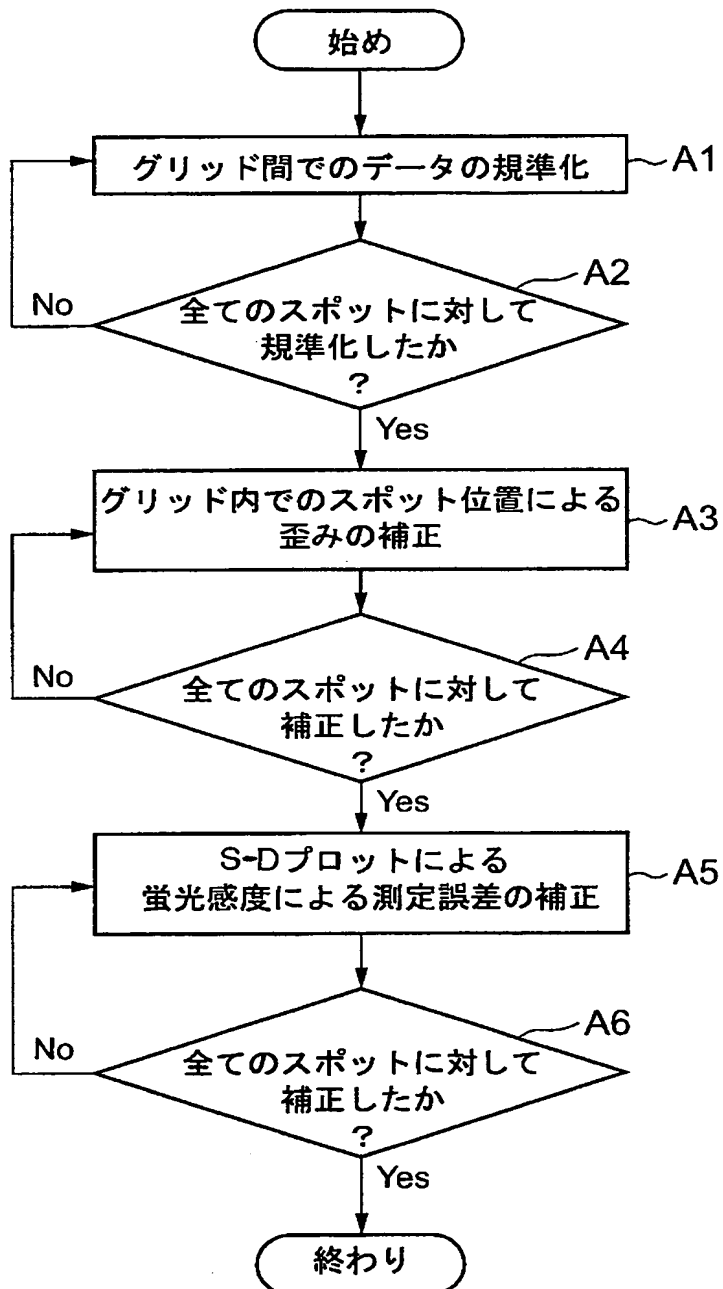




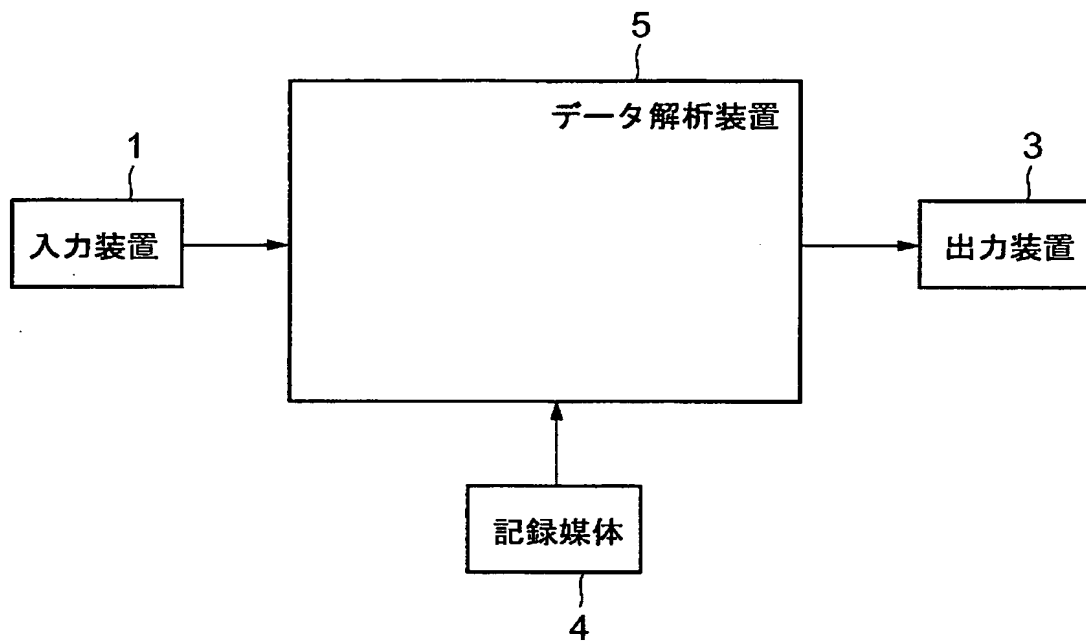
【図 2】



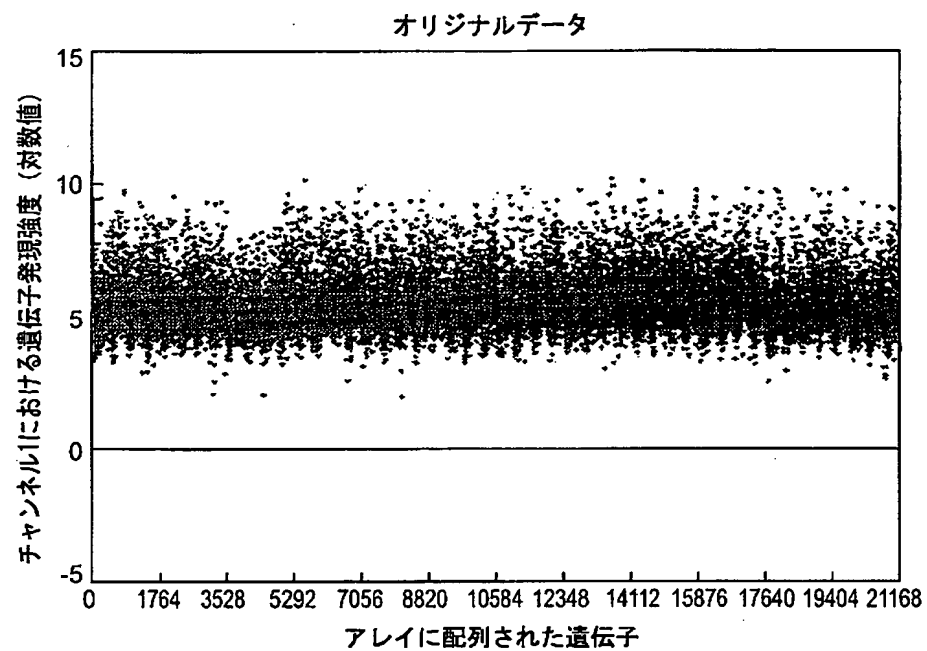
【図 3】



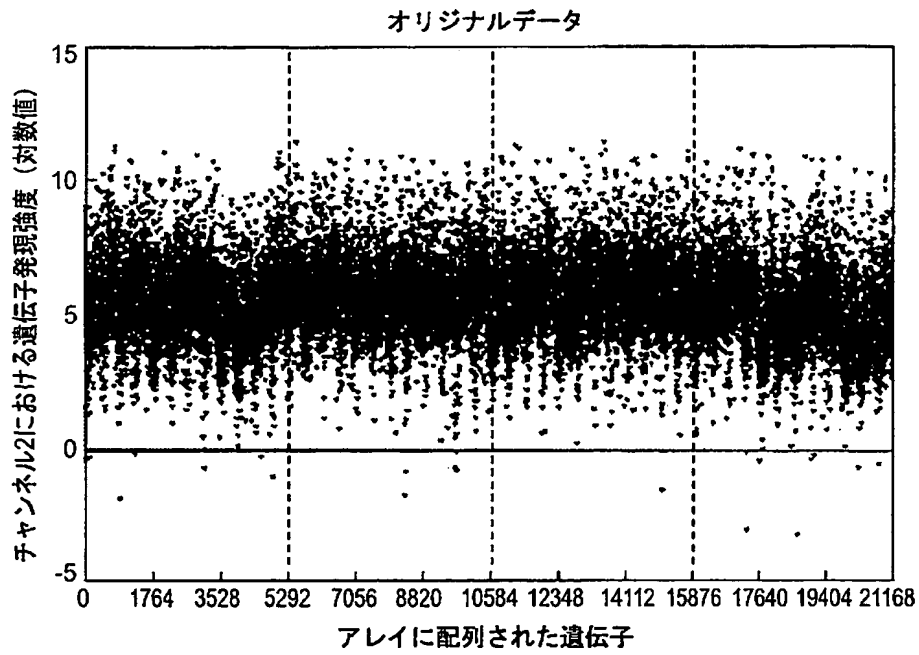
【図 4】



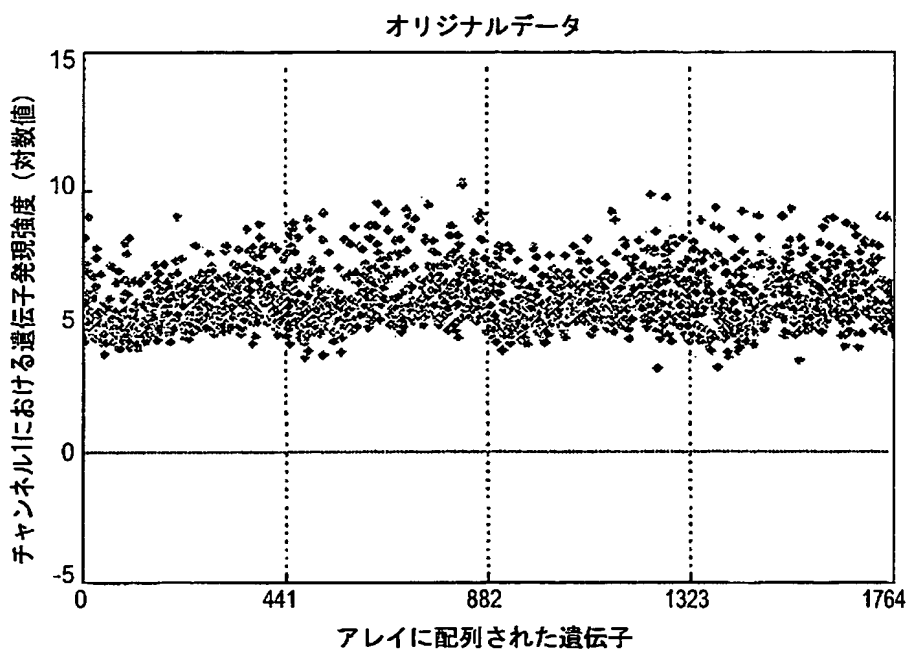
【図 5】



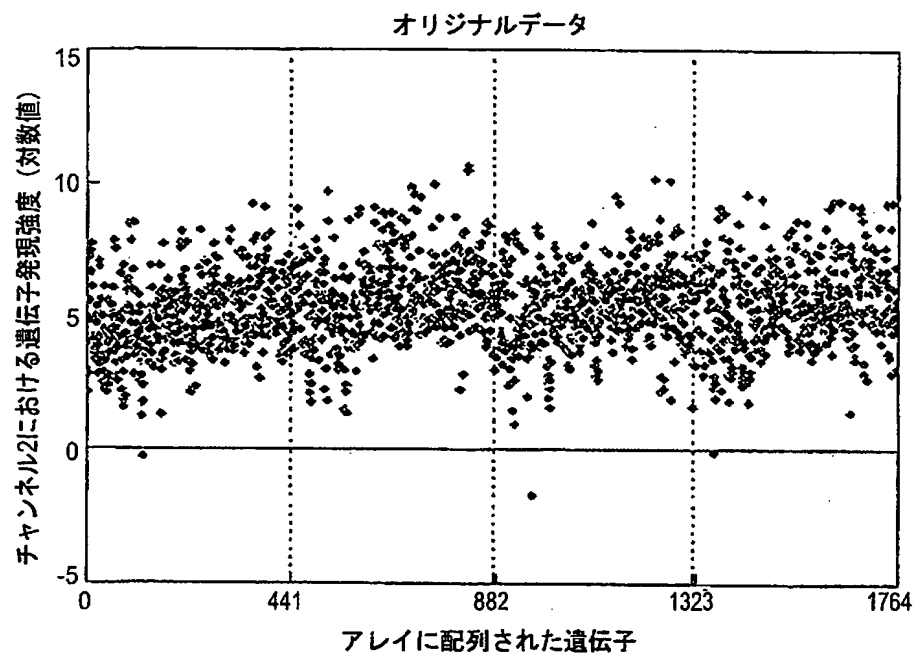
【図 6】



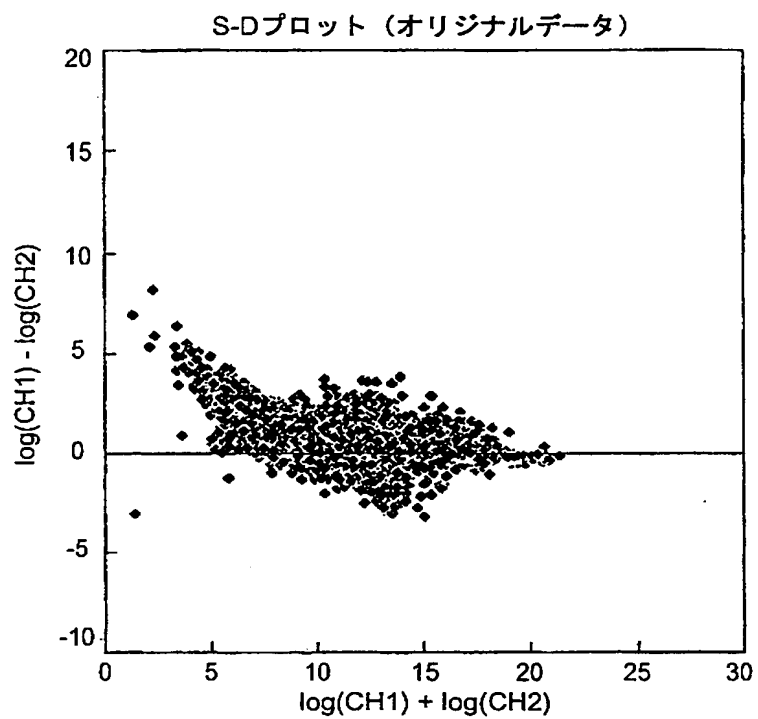
【図 7】



【図 8】

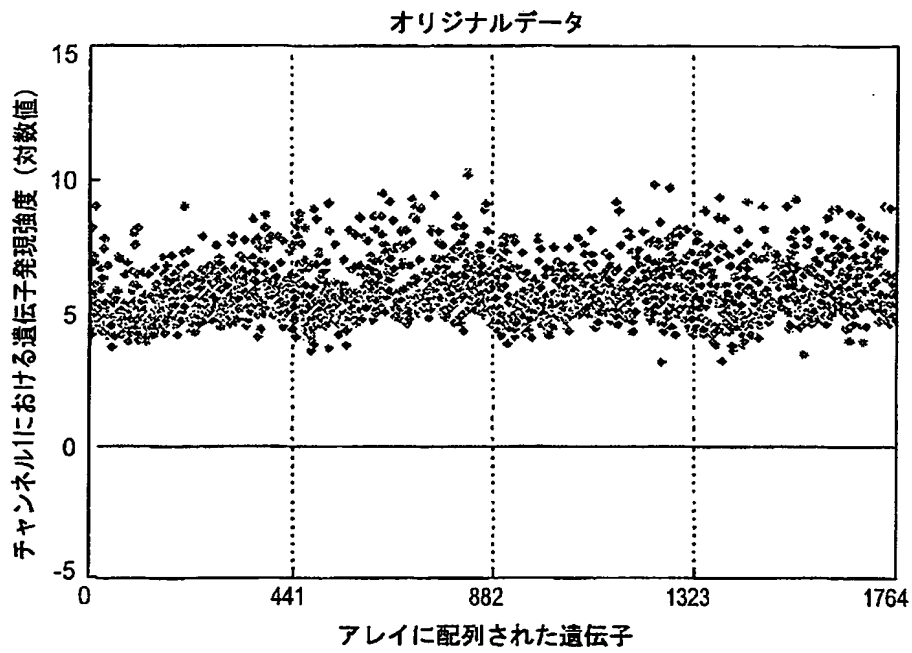


【図 9】

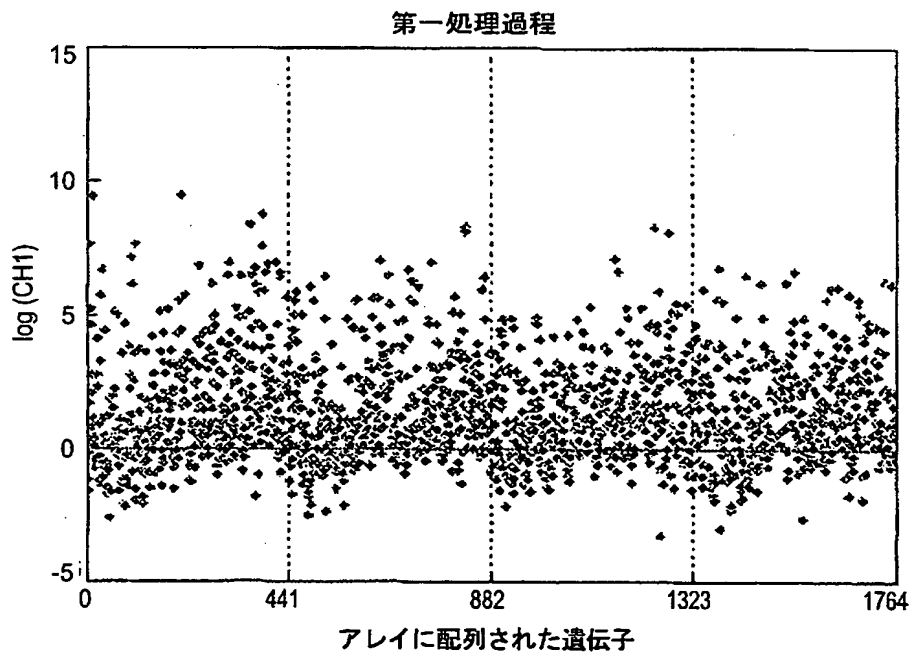




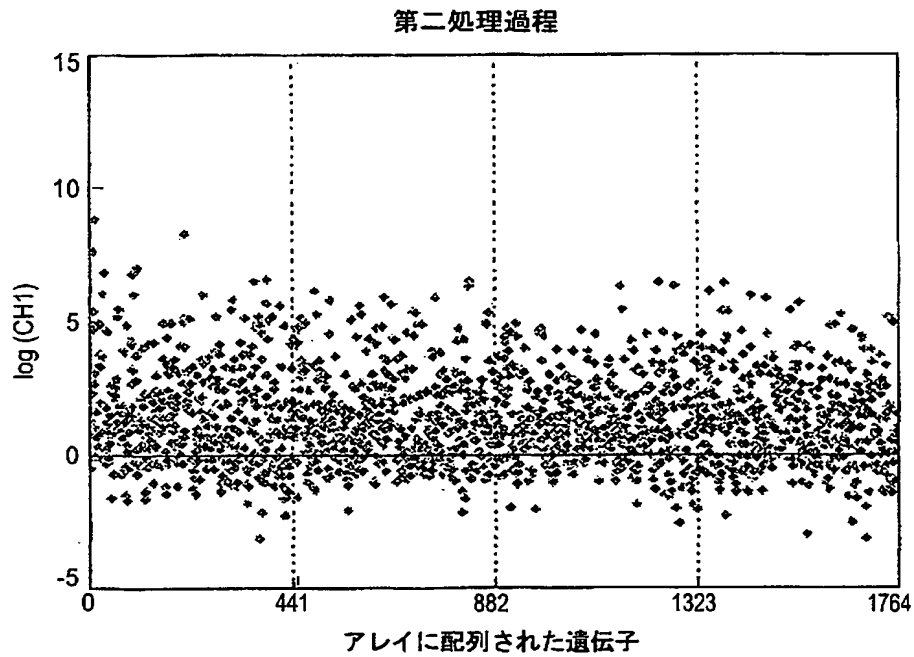
【図 10】



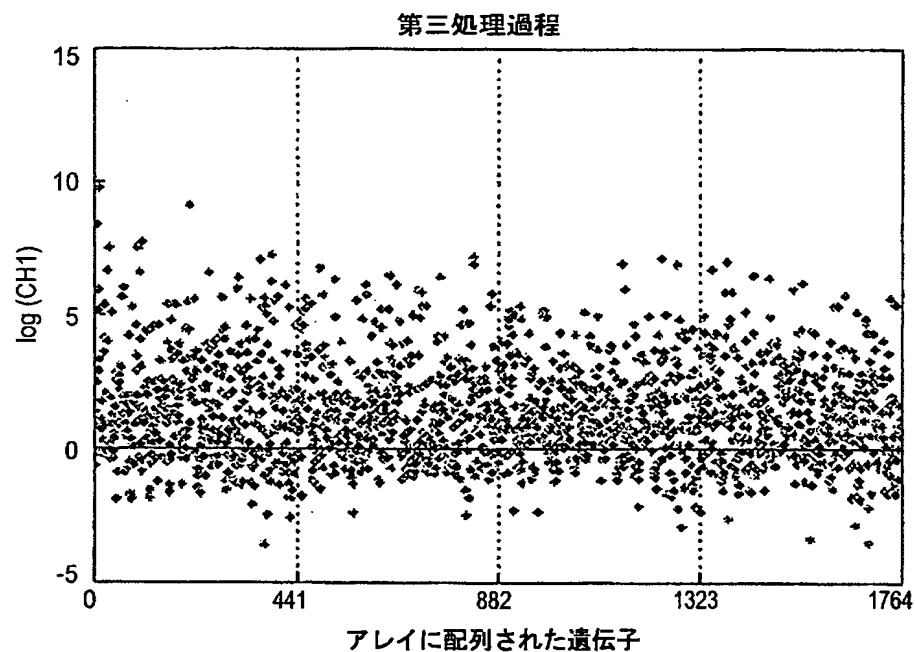
【図 11】



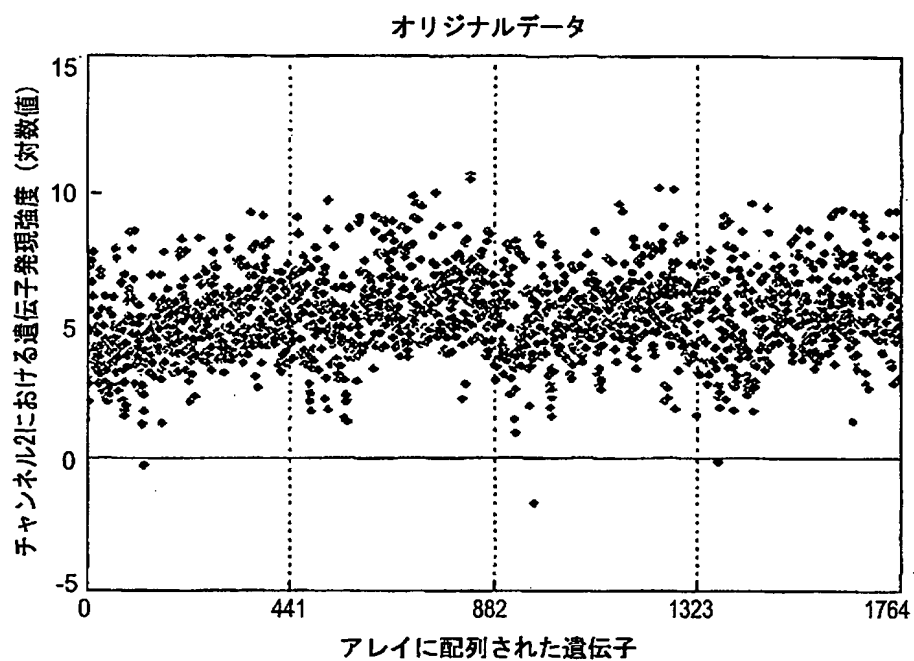
【図 12】



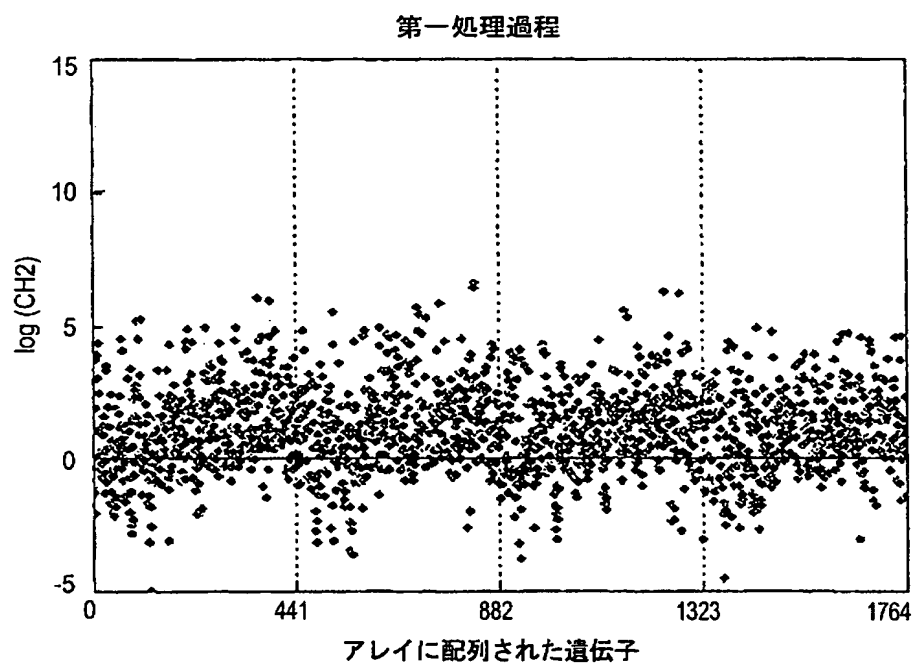
【図 13】



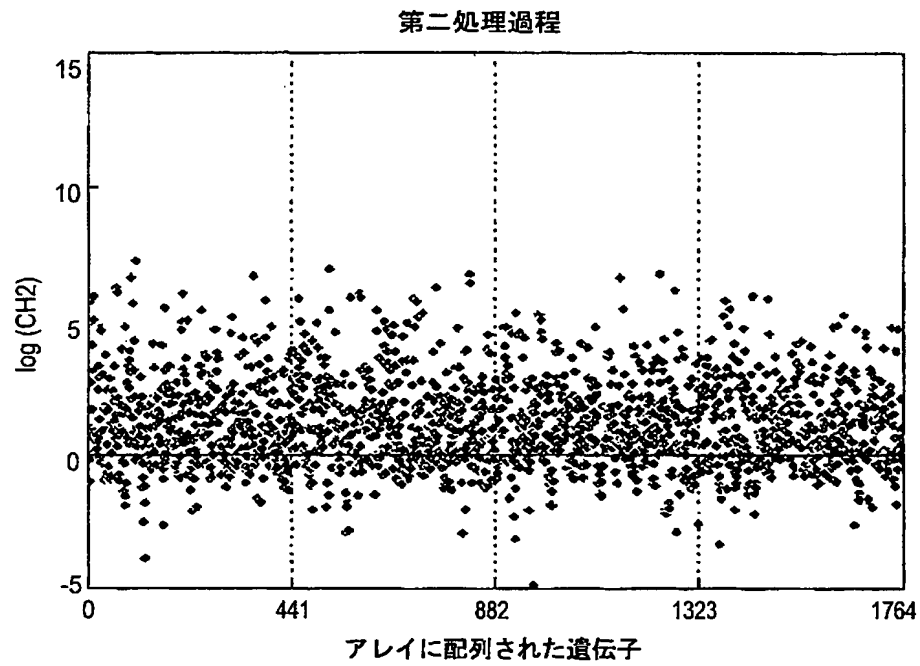
【図 14】



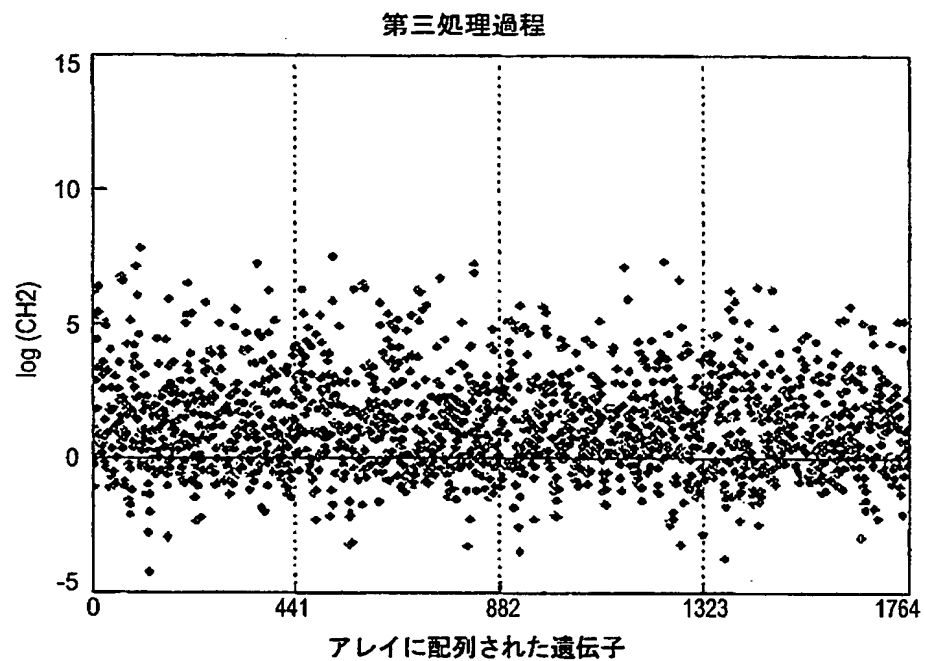
【図 15】



【図 16】

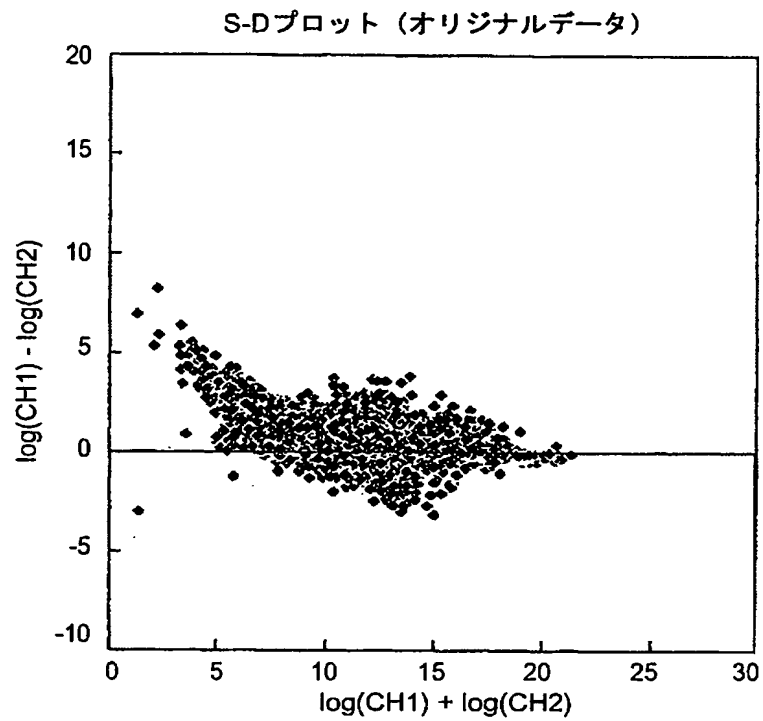


【図 17】

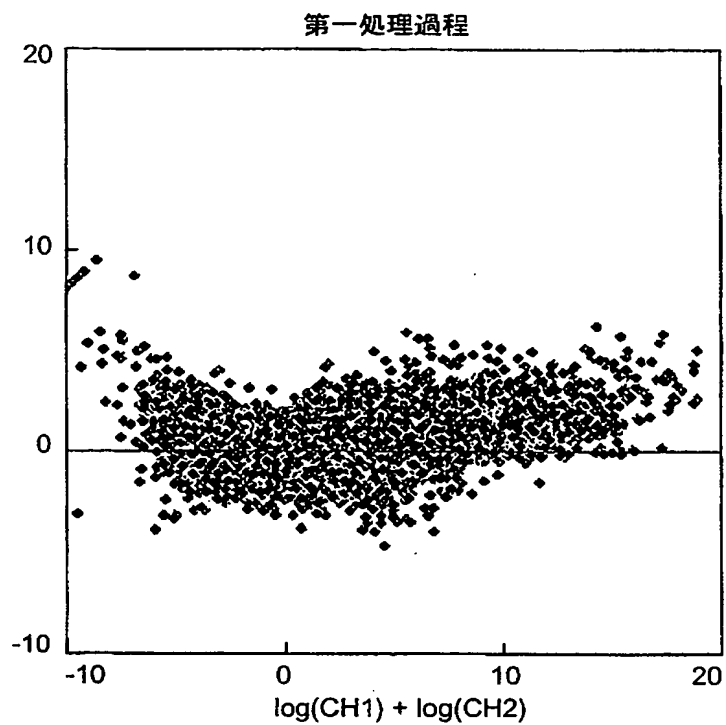




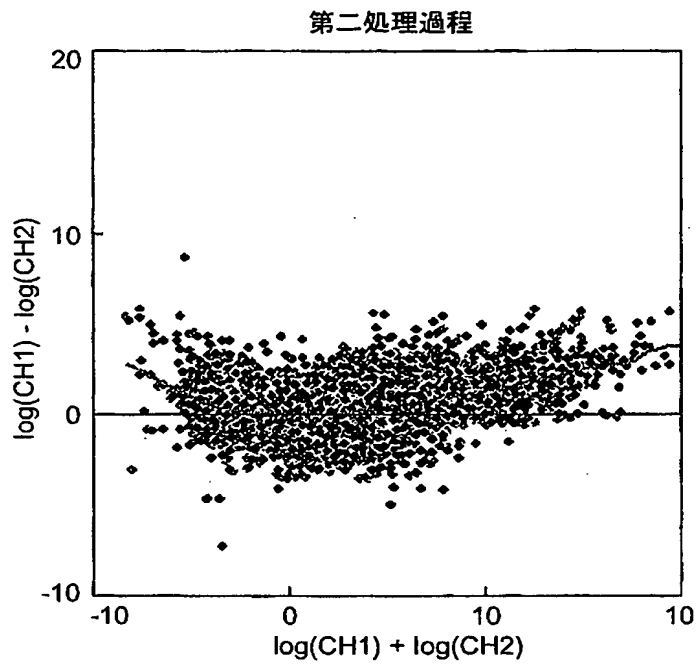
【図 18】



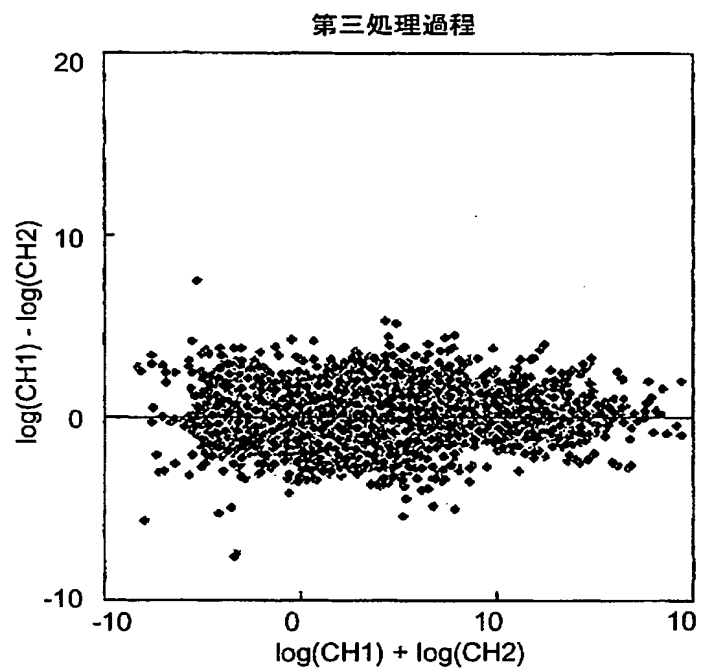
【図 19】



【図 2 0】



【図 2 1】



【書類名】 要約書

【要約】

【課題】 マイクロアレイデータのグローバルおよびローカルな歪みに対してより精密な補正をし、さらに蛍光色素の感度の違いによる測定誤差を補正する。

【解決手段】 第一処理過程であるデータ規準化手段は、入力装置から遺伝子発現強度データを入力し、大半の遺伝子は発現していないことを前提としてグリッド毎の順序統計量を用いて遺伝子発現強度データを規準化し、規準化した遺伝子発現強度データを出力する。第二処理過程であるスポット位置による補正手段は、グリッド毎にスポット位置に依存する歪みをノンパラメトリック平滑化法によって推定し、スポット位置に依存した歪みを補正した遺伝子発現強度データを出力する。第三処理過程である S-D プロットによる補正手段は、S-D 変換を行い、蛍光色素の感度の違いによる歪みをノンパラメトリック平滑化法によって推定し、蛍光色素の感度の違いによる歪みを補正した遺伝子発現強度データを出力装置に出力する。

【選択図】 図 2

特願 2 0 0 3 - 1 2 4 5 8 5

出 願 人 履 歴 情 報

識別番号

[ 0 0 0 0 0 4 2 3 7 ]

1. 変更年月日

1 9 9 0 年 8 月 2 9 日

[変更理由]

新規登録

住 所

東京都港区芝五丁目 7 番 1 号

氏 名

日本電気株式会社

特願 2 0 0 3 - 1 2 4 5 8 5

出 願 人 履 歴 情 報

識別番号

[ 5 0 3 0 7 7 1 6 5 ]

1. 変更年月日

2 0 0 3 年 2 月 2 6 日

[変更理由]

新規登録

住 所

広島県廿日市市宮園 9 丁目 1 の 7

氏 名

大瀧 慈

特願 2 0 0 3 - 1 2 4 5 8 5

出 願 人 履 歴 情 報

識別番号

[ 5 0 0 5 3 5 3 0 1 ]

1 . 変更年月日

2 0 0 0 年 1 1 月 2 0 日

[変更理由]

新規登録

住 所

東京都中央区八丁堀二丁目 2 6 番 9 号 グランデビルディング

氏 名

社団法人バイオ産業情報化コンソーシアム